

On Probability Estimation by Exponential Smoothing

Christopher Mattern

Technische Universität Ilmenau

Ilmenau, Germany

christopher.mattern@tu-ilmenau.de

Abstract. Probability estimation is essential for every statistical data compression algorithm. In practice probability estimation should be adaptive, i.e. recent observations should receive a higher weight than older observations. We present a probability estimation method based on exponential smoothing that satisfies this requirement and runs in constant time per letter. Our main contribution is a theoretical analysis in case of a binary alphabet for various smoothing rate sequences: We show that the redundancy w.r.t. a piecewise stationary model with s segments is $O(s\sqrt{n})$ for any bit sequence of length n , an improvement over redundancy $O(s\sqrt{n \log n})$ of previous approaches with similar time complexity.

1 Introduction

Background. Sequential probability assignment is an elementary component of every statistical data compression algorithm, such as Prediction by Partial Matching, Context Tree Weighting and PAQ (“pack”). Statistical compression algorithms split compression into modeling and coding and process an input sequence letter-by-letter. During modeling a model computes a distribution p and during coding an encoder maps the next letter x , given p , to a codeword of a length close to $-\log p(x)$ bits (this is the ideal code length). Decoding is the reverse: Given p and the codeword the decoder restores x . Arithmetic Coding is the de facto standard en-/decoder, it closely approximates the ideal code length [1]. All of the mentioned compression algorithms require simple, elementary models to predict a probability distribution. Elementary models are typically based on simple closed-form expressions, such as relative letter frequencies. Nevertheless, elementary models have a big impact on both theoretical guarantees [9, 11] and empirical performance [4, 10] of statistical compression algorithms. (Commonly, we express theoretical guarantees on a model by the amount of bits the model requires above an ideal competing scheme assuming ideal encoding, the so-called redundancy.) It is wise to choose elementary models carefully and desirable to analyze them theoretically and to study them experimentally. In this work we focus on elementary models with the ability to adapt to changing statistics (see next paragraph) whose implementation meets practical requirements, that is $O(Nn)$ (arithmetic) operations and $O(N)$ data words (holding e. g. a counter or a rational number) while processing a sequence of length n over an alphabet of size N .

Previous Work. Relative frequency-based elementary models, such as the Laplace- and KT-Estimator, are well-known and well-understood [1]. A major drawback of these classical techniques is that they don’t exploit recency effects (adaptivity): For an accurate prediction

novel observations are of higher importance than past observations [2, 7]. From a theoretical point of view adaptivity is evident in low redundancy w. r. t. an adaptive competing scheme such as a Piecewise Stationary Model (PWS). A PWS partitions a sequence of length n arbitrarily into s segments and predicts an arbitrary fixed probability distribution for every letter within a segment. (Since both segmentation and prediction within a segment are arbitrary, we may assume both to be optimal.)

To lift the limitation of classical relative frequency-based elementary models we typically age observation counts, aging takes place immediately before incrementing the count of a novel letter. For aging frequency-based elementary models there exist two major strategies, which are heavily used in practice. In *Strategy 1* (count rescaling) we divide all counts by a factor in well-defined intervals (e. g. when the sum of all counts exceeds a threshold) [2], for *Strategy 2* (count smoothing) we multiply all counts by a factor in $(0, 1)$ in every update [7]. Strategy 1 was analyzed in [5] and has redundancy $O(s\sqrt{n} \log n)$. Similarly, a KT-estimator which completely discards all counts periodically was analyzed in [8] and has redundancy $O(s\sqrt{n} \log n)$. Strategy 2 was studied mainly experimentally [7, 9, 10].

Another approach for adaptive probability estimation, commonly used in PAQ, is smoothing of probabilities, *Strategy 3*. Given a probability distribution (i. e. the prediction of the previous step) and a novel letter we carry out an update as follows: First we multiply all probabilities with smoothing rate $\alpha \in (0, 1)$ and afterwards we increment the probability of the novel letter by $1 - \alpha$. Smoothing rate α does not vary from step to step. To our knowledge this common-sense approach was first mentioned in [3]. A finite state machine that approximates smoothing was analyzed in [6] and has redundancy $O(nK^{-2/3})$ w. r. t. PWS with $s = 1$, where K is the number of states.

All aforementioned approaches meet practical demands, they require $O(Nn)$ (arithmetic) operations and $O(N)$ data words. More complex (but unpractical) methods are based on mixtures over elementary models associated to so-called transition diagrams [8, 12] or associated to (PWS-)partitions [10].

Our Contribution. In this work we analyze a generalization of strategies 2 and 3 for a binary alphabet. Based on mild assumptions on sequence $\alpha_1, \alpha_2, \dots$ of smoothing rates (α_k is used for an update after observing the k -th letter) we explicitly identify an input sequence with maximum redundancy w. r. t. PWS with $s = 1$ and subsequently derive redundancy bounds for $s \geq 1$ (Section 3). For PWS with arbitrary s we give redundancy bounds for three choices of smoothing rates in Section 4. First, we consider a fixed smoothing rate $\alpha = \alpha_1 = \alpha_2 = \dots$ (as in PAQ) and provide $\alpha^*(n)$ that guarantees redundancy $O(s\sqrt{n})$ for a sequence of length n ; second, we propose a varying smoothing rate, where $\alpha_k \approx \alpha^*(k)$; and finally a varying smoothing rate that is equivalent to Strategy 2 from the previous section. By tuning parameters we obtain redundancy $O(s\sqrt{n})$ for all smoothing rate choices, an improvement over redundancy guarantees known so far for models requiring $O(Nn)$ (arithmetic) operations per input sequence. Section 5 supports our bounds with a small experimental study and finally Section 6 summarizes and evaluates our results and gives perspectives for future research.

2 Preliminaries

Sequences. We use $x_{i:j}$ to denote a sequence $x_i x_{i+1} \dots x_j$ of objects (numbers, letters, \dots). Unless stated differently, sequences are bit sequences (have letters $\{0, 1\}$). If $i > j$, then

$x_{i:j} := \phi$, where ϕ is the empty sequence; if $j = \infty$, then $x_{i:j} = x_i x_{i+1} \dots$ has infinite length. For sequence $x_{1:n}$ define $x_{<i} := x_{1:i-1}$ and $x_{\leq i} := x_{1:i}$; we call $x_{1:n}$ *deterministic*, if $x_1 = \dots = x_n$, and *non-deterministic*, otherwise.

Code Length and Entropy. Code length is measured in bits, thus $\log := \log_2$. For probability distribution p over $\{0, 1\}$ and letter x we define $\ell(x; p) := -\log p(x)$. The binary entropy function is denoted as $H(q) := -q \log(q) - (1-q) \log(1-q)$, for a probability q . For sequence $x_{1:n}$ and relative frequency q of a 1-bit in $x_{1:n}$ let $h(x_{1:n}) := nH(q)$ be the *empirical entropy* of $x_{1:n}$.

Partitions and Sets. Calligraphic letters denote sets. A partition of a non-empty segment (interval) $(a, b]$ of integers is a set of non-overlapping segments $(i_0, i_1], \dots, (i_{n-1}, i_n]$ s. t. $a = i_0 < i_1 < \dots < i_n = b$. The phrase *k-th segment* uniquely refers to $(i_{k-1}, i_k]$.

Models and Exponential Smoothing of Probabilities. We first characterize the term model from statistical data compression, in order to define our modeling method. A *model* MDL maps a sequence $x_{1:n}$ of length $n \geq 0$ to a probability distribution p on $\{0, 1\}$. We define the short-hands $\text{MDL}(x_{\leq n}) := p$ (this is **not** the probability of sequence $x_{\leq n}$!) and $\text{MDL}(x; x_{\leq n}) := p(x)$. Model MDL assigns $\ell(x_{\leq n}; \text{MDL}) := -\sum_{1 \leq k \leq n} \log \text{MDL}(x_k; x_{<k})$ bits to sequence $x_{\leq n}$. We are now ready to formally define our model of interest.

Definition 2.1. For sequence $\alpha_{1:\infty}$, where $0 < \alpha_1, \alpha_2, \dots < 1$, and probability distribution p , where $p(0), p(1) > 0$, we define model $\text{ESP} = (\alpha_{1:\infty}, p)$ by the sequential probability assignment rule

$$\text{ESP}(x; x_{\leq k}) = \begin{cases} \alpha_k \text{ESP}(x; x_{<k}) + 1 - \alpha_k, & \text{if } k > 0 \text{ and } x = x_k \\ \alpha_k \text{ESP}(x; x_{<k}), & \text{if } k > 0 \text{ and } x \neq x_k \\ p(x), & \text{if } k = 0 \end{cases} \quad (1)$$

Smoothing rates control the adaption of ESP, large α_i 's give high weight to old observations and low weight to new observations, the converse holds for small α_i 's. For our analysis we must assume that the smoothing rates are sufficiently large:

Assumption 2.2. $\text{ESP} = (\alpha_{1:\infty}, p)$ satisfies $\frac{1}{2} < \alpha_1, \alpha_2, \dots < 1$ and w. l. o. g. $p(0) \leq p(1)$.

For the upcoming analysis the product of smoothing rates plays an important role. Hence, given smoothing rate sequence $\alpha_{1:\infty}$, we define $\beta_0 = 1$ and $\beta_i := \alpha_1 \cdot \dots \cdot \alpha_i$, for $i > 0$.

3 Redundancy Analysis

First Main Result. Now we can state our first main result which compares ESP to the code length of an optimal fixed code for $x_{1:n}$, that is the empirical entropy $h(x_{1:n})$. Before we prove the theorem, we discuss its implications.

Theorem 3.1. If Assumption 2.2 holds, then we have

$$\begin{aligned} & \ell(x_{1:n}; \text{ESP}) - h(x_{1:n}) \\ & \leq \begin{cases} \sum_{i=0}^{n-1} \log \frac{1}{1-p(1)\beta_i}, & \text{if } x_{1:n} \text{ is deterministic} \\ \log \frac{1}{p(1)\beta_{n-1}} + \sum_{i=0}^{n-2} \log \frac{1}{1-p(1)\beta_i} - nH\left(\frac{1}{n}\right), & \text{otherwise} \end{cases} \end{aligned} \quad (2)$$

Recall that by Assumption 2.2 we have $p(0) \leq p(1)$. First, consider a deterministic sequence of 0-bits. By (1) we have $\text{ESP}(1; x_{\leq i}) = \beta_i p(1)$, thus $\text{ESP}(0; x_{\leq i}) = 1 - \beta_i p(1)$, both for $0 \leq i < n$. The total code length is $\ell(x_{1:n}; \text{ESP}) = -\sum_{i=0}^{n-1} \log(1 - \beta_i p(1))$ and clearly $h(x_{1:n}) = 0$, so the redundancy $\ell(x_{1:n}; \text{ESP}) - h(x_{1:n})$ matches (2). Now consider a non-deterministic sequence $x_{1:n} = 00 \dots 01$ with single 1-bit at position n . Similar to the deterministic case the same equations for $\text{ESP}(\cdot; x_{\leq i})$ hold, for $i < n$. The total code length is $\ell(x_{1:n}; \text{ESP}) = -\sum_{i=0}^{n-2} \log(1 - \beta_i p(1)) - \log(\beta_{n-1} p(1))$ and the empirical entropy is $h(x_{1:n}) = nH(\frac{1}{n})$. Again, the redundancy matches (2). In summary, if $p(0) \leq p(1)$, then $00 \dots 0$ is a deterministic sequence with maximum redundancy and $00 \dots 01$ is a non-deterministic sequence with maximum redundancy. Similar statements hold, if $p(0) \geq p(1)$: by symmetry we must toggle 0-bits and 1-bits. When $p(0) = p(1)$ we have equal redundancy, e. g. $\ell(00 \dots 0; \text{ESP}) = \ell(11 \dots 1; \text{ESP})$ in the deterministic case. In summary, for a given instance of ESP (that satisfies Assumption 2.2) the worst-case input is either $00 \dots 0$ (only 0-bits) or $00 \dots 01$ (single 1-bit), among all 2^n bit sequences of length n . For fixed n we can now easily compare the redundancies of those two inputs and immediately depict the worst-case input and its redundancy.

For the proof of Theorem 3.1 we require the following lemma.

Lemma 3.2. *Any non-deterministic sequence $x_{1:n}$ of length $n \geq 2$ satisfies*

$$h(x_{1:n}) - h(x_{2:n}) \geq \begin{cases} nH\left(\frac{1}{n}\right), & \text{if } x_{2:n} \text{ is deterministic} \\ nH\left(\frac{1}{n}\right) - (n-1)H\left(\frac{1}{n-1}\right), & \text{otherwise} \end{cases}.$$

Proof. Let $1-p$ be the relative frequency of x_1 in $x_{1:n}$, thus $h(x_{1:n}) - h(x_{2:n}) = nH(p) - (n-1)H\left(\frac{n}{n-1} \cdot p\right) =: f(p)$. We distinguish two cases:

Case 1: $x_{2:n}$ is deterministic. We have $p = \frac{n-1}{n}$ and $f(p) = nH\left(\frac{n-1}{n}\right) = nH\left(\frac{1}{n}\right)$.

Case 2: $x_{2:n}$ is non-deterministic. Since $H(p)$ is concave, $H'(p)$ is decreasing and $f'(p) = n[H'(p) - H'\left(\frac{n}{n-1} \cdot p\right)] \geq 0$, i. e. $f(p)$ is increasing and minimal for minimum p . Since $x_{1:n}$ is non-deterministic the minimum value of p is $\frac{1}{n}$ and we get $f(p) \geq f\left(\frac{1}{n}\right) = nH\left(\frac{1}{n}\right) - (n-1)H\left(\frac{1}{n-1}\right)$. \square

Now let us proceed with the major piece of work in this section.

Proof of Theorem 3.1. We define $r(x_{1:n}, \text{ESP}) := \ell(x_{1:n}; \text{ESP}) - h(x_{1:n})$ and distinguish:

Case 1: $x_{1:n}$ is deterministic. By $p(0) \leq p(1)$ (Assumption 2.2) we have $\text{ESP}(x; x_{< i}) \geq \text{ESP}(0; x_{< i}) = 1 - p(1)\beta_{i-1}$ and $h(x_{1:n}) = 0$, we get

$$r(x_{1:n}, \text{ESP}) = \sum_{1 \leq i \leq n} \log \frac{1}{\text{ESP}(x_i; x_{< i})} \leq \sum_{0 \leq i < n} \log \frac{1}{1 - p(1)\beta_i}.$$

Case 2: $x_{1:n}$ is non-deterministic. We have $n \geq 2$ and by induction on n we prove

$$r(x_{1:n}, \text{ESP}) \leq \log \frac{1}{p(1)\beta_{n-1}} + \sum_{0 \leq i < n-1} \log \frac{1}{1 - p(1)\beta_i} - nH\left(\frac{1}{n}\right).$$

Base: $n = 2$. — We have $x_{1:n} \in \{01, 10\}$, in either case $h(x_{1:n}) = nH\left(\frac{1}{n}\right) = 2$ and $\ell(x_{1:n}; \text{ESP}) = \log \frac{1}{p(x_1)\beta_1 p(x_2)} = \log \frac{1}{p(1)\beta_1} + \log \frac{1}{1-p(1)}$, the claim follows.

Step: $n > 2$. — By defining $\text{ESP}' = (\alpha'_{1:\infty}, p')$, where $\alpha'_i = \alpha_{i+1}$, $\beta'_i = \alpha'_1 \cdot \dots \cdot \alpha'_i$, $p' = \text{ESP}(x_{\leq 1})$ we may write

$$r(x_{1:n}, \text{ESP}) = \log \frac{1}{p(x_1)} + r(x_{2:n}, \text{ESP}') - (h(x_{1:n}) - h(x_{2:n})). \quad (3)$$

Now w.l.o.g. fix p' s.t. $p'(0) \leq p'(1)$. Since we want to bound (3) from above, we must choose x_1 s.t. $p(x_1)$ is minimal (and the r.h.s. of (3) is maximal). To do so, distinguish:

Case 1: $x_1 = 0$. For some distribution q with $q(0) > 0$ we have $p(x_1) = q(0)$ and $\frac{1}{2} \geq p'(0) = \alpha_1 q(0) + 1 - \alpha_1$, thus $q(0) \leq [\alpha_1 - \frac{1}{2}] / \alpha_1$. (Notice the subtle detail: $\alpha_1 \leq \frac{1}{2}$ implies $q(0) \leq 0$, which contradicts $q(0) > 0$ and would make Case 1 impossible; however we assumed $\alpha_1 > \frac{1}{2}$.) Furthermore, we have $q(0) \leq \frac{1}{2}$.

Case 2: $x_1 = 1$. For some distribution r with $r(1) > 0$ we have $p(x_1) = r(1)$ and $\frac{1}{2} \leq p'(1) = \alpha_1 r(1) + 1 - \alpha_1$, thus $r(1) \geq [\alpha_1 - \frac{1}{2}] / \alpha_1$.

Since $q(0) \leq r(1)$ (i.e. Case 1 minimizes $p(x_1)$) and $q(0) \leq \frac{1}{2}$ we may now w.l.o.g. assume that $x_1 = 0, p'(1) = \alpha_1 p(1), p(x_1) = 1 - p(1)$ and $p(0) \leq p(1)$. We distinguish:

Case 1: $x_{2:n}$ is deterministic. We must have $x_{2:n} = 11 \dots 1$, since $x_1 = 0$ and $x_{1:n}$ is non-deterministic, thus

$$r(x_{2:n}, \text{ESP}') = \sum_{0 \leq i < n-1} \log \frac{1}{1 - p'(0)\beta'_i} \leq \log \frac{1}{p(1)\beta_{n-1}} + \sum_{1 \leq i < n-1} \log \frac{1}{1 - p(1)\beta_i}, \quad (4)$$

where we obtain the inequality by $p'(0)\beta'_i \leq p'(1)\beta'_i = p(1)\beta_{i+1}$, for $i < n-2$ and $1 - p'(0)\beta'_{n-2} = 1 - [1 - p(1)\alpha_1]\beta_{n-1}/\alpha_1 \geq p(1)\beta_{n-1}$, for $i = n-2$. To obtain the claim we plug the inequalities (4) and $h(x_{1:n}) - h(x_{2:n}) \geq nH(\frac{1}{n})$ (by Lemma 3.2) into (3) and note that $p(x_1) = 1 - p(1)\beta_0$ (since $\beta_0 = 1$).

Case 2: $x_{2:n}$ is non-deterministic. The hypothesis and $p'(1)\beta'_i = p(1)\beta_{i+1}$ yield

$$\begin{aligned} r(x_{2:n}, \text{ESP}') &\leq \log \frac{1}{p'(1)\beta'_{n-2}} + \sum_{0 \leq i < n-2} \log \frac{1}{1 - p'(1)\beta'_i} - (n-1)H\left(\frac{1}{n-1}\right) \\ &= \log \frac{1}{p(1)\beta_{n-1}} + \sum_{1 \leq i < n-1} \log \frac{1}{1 - p(1)\beta_i} - (n-1)H\left(\frac{1}{n-1}\right). \end{aligned} \quad (5)$$

We plug the inequalities (5) and $h(x_{1:n}) - h(x_{2:n}) \geq nH(\frac{1}{n}) - (n-1)H(\frac{1}{n-1})$ (by Lemma 3.2), into (3) and note that $p(x_1) = 1 - p(1)\beta_0$ (since $\beta_0 = 1$) to end the proof. \square

Second Main Result. Let us now extend the competing scheme of Theorem 3.1, to which we compare ESP to. Suppose the competing scheme splits the input sequence $x_{1:n}$ according to an arbitrary partition \mathcal{S} of $[1, n]$ and may use an optimal fixed code within every segment $[a, b] \in \mathcal{S}$. The competing scheme has total coding cost $h(x_{a:b})$ for $x_{a:b}$, thus coding cost $\sum_{[a,b] \in \mathcal{S}} h(x_{a:b})$ for $x_{1:n}$. Notice, that this is a lower bound on the coding cost of any PWS with partition \mathcal{S} . Since the situation within a segment resembles the situation of Theorem 3.1, we may now naturally extend the redundancy analysis to the aforementioned competitor.

Theorem 3.3. *Let \mathcal{S} be an arbitrary partition of $[1, n]$. If Assumption 2.2 holds, then*

$$\ell(x_{1:n}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:b}) \leq |\mathcal{S}| \log \frac{1}{p(0)\beta_{n-1}} + \sum_{(a,b] \in \mathcal{S}} \sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a}. \quad (6)$$

Proof. Let $r(x_{1:n}, \text{ESP}) := \ell(x_{1:n}; \text{ESP}) - h(x_{1:n})$. Our plan for the proof is to simplify (2) (see calculations below) to yield

$$\ell(x_{1:n}; \text{ESP}) - h(x_{1:n}) \leq \log \frac{1}{p(0)\beta_{n-1}} + \sum_{1 \leq i < n} \log \frac{1}{1 - \beta_i} \quad (7)$$

and use to (7) to bound the redundancy for an arbitrary segment $(a, b]$ from \mathcal{S} (see calculations below) via

$$\sum_{a < i \leq b} \log \frac{1}{\text{ESP}(x_i; x_{<i})} - h(x_{a+1:b}) \leq \log \frac{1}{p(0)\beta_{b-1}} + \sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a}. \quad (8)$$

We now obtain (6) easily by summing (8) over all segments $(a, b]$ from \mathcal{S} and by $\beta_{b-1} \geq \beta_{n-1}$.

Simplifying (2). Observe that $\sum_{0 \leq i < n} \log \frac{1}{1 - p(1)\beta_i} = \log \frac{1}{p(0)} + \sum_{1 \leq i < n} \log \frac{1}{1 - p(1)\beta_i}$ and furthermore $p(0) \leq p(1)$. So bound (2) becomes

$$r(x_{1:n}, \text{ESP}) \leq \log \frac{1}{p(0)} + \sum_{1 \leq i < n} \log \frac{1}{1 - p(1)\beta_i} \leq \log \frac{1}{p(0)} + \sum_{1 \leq i < n} \log \frac{1}{1 - \beta_i},$$

if $x_{1:n}$ is deterministic and by $\log \frac{1}{p(1)} - nH\left(\frac{1}{n}\right) \leq 0$ (since $p(1) \geq \frac{1}{2}$ and $n \geq 2$)

$$\begin{aligned} r(x_{1:n}, \text{ESP}) &\leq \log \frac{1}{p(0)p(1)\beta_{n-1}} + \sum_{1 \leq i < n-1} \log \frac{1}{1 - p(1)\beta_i} - nH\left(\frac{1}{n}\right) \\ &\leq \log \frac{1}{p(0)\beta_{n-1}} + \sum_{1 \leq i < n} \log \frac{1}{1 - \beta_i}, \end{aligned}$$

if $x_{1:n}$ is non-deterministic. In either case bound (7) holds.

Redundancy of $(a, b]$. For segment $(a, b]$ we define sequence $x'_{1:b-a} = x_{a+1:b}$ and $\text{ESP}' = (\alpha'_{1:\infty}, p')$, s. t. $\text{ESP}(x; x_{<i}) = \text{ESP}'(x'; x'_{<i-a})$ for $i \in (a, b]$. Therefore, let $\alpha'_{1:\infty} = \alpha_{a+1:\infty}$, $\beta'_i = \alpha'_1 \cdot \dots \cdot \alpha'_i$ and w. l. o. g. $p'(0) \leq p'(1)$. We obtain

$$\begin{aligned} &\sum_{a < i \leq b} \log \frac{1}{\text{ESP}(x_i; x_{<i})} - h(x_{a+1:b}) \\ &= \sum_{1 \leq i \leq b-a} \log \frac{1}{\text{ESP}'(x'_i; x'_{<i})} - h(x'_{1:b-a}) = \ell(x_{1:n}, \text{ESP}) - h(x_{1:n}) \\ &\stackrel{(7)}{\leq} \log \frac{1}{p'(0)\beta'_{b-a-1}} + \sum_{1 \leq i < b-a} \log \frac{1}{1 - \beta'_i} \leq \log \frac{1}{p(0)\beta_{b-1}} + \sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a}, \end{aligned}$$

where the last step is due to $p'(0) \geq p(0)\beta_a$ (also $p'(1) \geq p(0)\beta_a$) and $\beta'_i = \beta_{a+i}/\beta_a$. \square

4 Choice of Smoothing Rate Sequence

Fixed Smoothing Rate. A straight-forward choice for the smoothing rates is to use the same rate α in every step. This leads to a simple and fast implementation, since no smoothing rate sequence needs to be computed or stored. We require the following lemma for the analysis:

Lemma 4.1. For $0 < \alpha < 1$ we have $\sum_{1 \leq i \leq m} \log \frac{1}{1 - \alpha^i} \leq \frac{(\pi \log e)^2}{6 \log \frac{1}{\alpha}}$.

Proof. For $m = 0$ the bound trivially holds, let $m \geq 1$. Since $\log \frac{1}{1-\alpha^z}$ is decreasing in z and integrable for z in $[0, \infty)$ we may bound the series by an integral,

$$\sum_{1 \leq i \leq m} \log \frac{1}{1-\alpha^i} \leq \int_0^m \log \frac{1}{1-\alpha^z} dz = \log(e) \int_0^m \sum_{j \geq 1} \frac{\alpha^{jz}}{j} dz. \quad (9)$$

The equality in (9) follows from the series expansion $\ln \frac{1}{1-y} = \sum_{j \geq 1} y^j/j$, for $|y| < 1$. To end the proof, it remains to bound the integral in (9) as follows (notice $\sum_{j \geq 1} j^{-2} = \pi^2/6$):

$$\int_0^m \sum_{j \geq 1} \frac{\alpha^{jz}}{j} dz = \sum_{j \geq 1} \frac{1}{j} \int_0^m \alpha^{jz} dz = \frac{\log e}{\log \frac{1}{\alpha}} \sum_{j \geq 1} \frac{1 - \alpha^{jm}}{j^2} \leq \frac{\pi^2 \log e}{6 \log \frac{1}{\alpha}}. \quad \square$$

Corollary 4.2. *Let \mathcal{S} be an arbitrary partition of $[1, n]$. If $\alpha = \alpha_1 = \alpha_2 = \dots$ and Assumption 2.2 holds, then*

$$\ell(x_{1:n}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:b}) \leq |\mathcal{S}| \cdot \left[\log \frac{1}{p(0)} + \frac{(\pi \log e)^2}{6 \log \frac{1}{\alpha}} + (n-1) \log \frac{1}{\alpha} \right]. \quad (10)$$

Proof. We have $\beta_i = \alpha^i$, thus for $i \in (a, b]$ we plug the estimate

$$\sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a} = \sum_{0 < i-a < b-a} \log \frac{1}{1 - \alpha^{i-a}} \stackrel{\text{Lem. 4.1}}{\leq} \frac{(\pi \log e)^2}{6 \log \frac{1}{\alpha}}$$

and $\log \beta_{n-1} = (n-1) \log \alpha$ into (6) to conclude the proof. \square

Choosing $\alpha = e^{-\frac{\pi}{\sqrt{6(n-1)}}}$ minimizes the r. h. s. of bound (10) and satisfies $\alpha > \frac{1}{2}$ (Assumption 2.2), when $n \geq 5$. The optimal choice gives redundancy at most

$$|\mathcal{S}| \cdot \left[\frac{2\pi \log e}{\sqrt{6}} \cdot \sqrt{n} + \log \frac{1}{p(0)} \right] < |\mathcal{S}| \cdot \left[3.701 \cdot \sqrt{n} + \log \frac{1}{p(0)} \right]. \quad (11)$$

Varying Smoothing Rate. It is impossible to choose an optimal fixed smoothing rate, when n is unknown. A standard technique to handle this situation is the doubling trick, which will increase the \sqrt{n} -term in (11) by a factor of $\sqrt{2}/(\sqrt{2}-1) \approx 3.41$. However, we can do better by slowly increasing the smoothing rate step-by-step, which only leads to a factor $\sqrt{2} \approx 1.41$.

Corollary 4.3. *Let \mathcal{S} be an arbitrary partition of $[1, n]$. If $\alpha_k = e^{-\pi/\sqrt{12(k+1)}}$ (i. e. $\alpha_k > \frac{1}{2}$) and Assumption 2.2 holds, then*

$$\ell(x_{1:n}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:b}) \leq |\mathcal{S}| \cdot \left[\log \frac{1}{p(0)} + \frac{2\pi \log e}{\sqrt{3}} \cdot \sqrt{n} \right]. \quad (12)$$

Proof. We have $\beta_i = \exp\left(-\frac{\pi}{\sqrt{12}} \sum_{1 < k \leq i+1} k^{-1/2}\right)$ and bound the terms depending on the β_i 's in (6) from above. First, observe that

$$\sum_{1 < k \leq n} k^{-1/2} \leq \int_1^n \frac{dz}{\sqrt{z}} \leq 2\sqrt{n} \stackrel{\text{Def. } \beta_i}{\Rightarrow} \log \frac{1}{\beta_{n-1}} \leq \frac{\pi \log e}{\sqrt{3}} \sqrt{n}, \quad (13)$$

second, for $a < i < b$ we have $\beta_i/\beta_a = \alpha_{a+1} \cdot \dots \cdot \alpha_i \leq (\alpha_{n-1})^{i-a}$, since $i < n$ and $\alpha_1, \alpha_2, \dots$ is increasing, consequently we obtain

$$\sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a} \leq \sum_{a < i < b} \log \frac{1}{1 - (\alpha_{n-1})^{i-a}} \stackrel{\text{Lem. 4.1}}{\leq} \frac{(\pi \log e)^2}{6 \log \frac{1}{\alpha_{n-1}}} = \frac{\pi \log e}{\sqrt{3}} \sqrt{n}. \quad (14)$$

We plug (13) and (14) into (6), the result is (12). \square

Count Smoothing. Consider aging Strategy 2 from Section 1 with smoothing rate $\lambda \in (0, 1)$. We will now show that Strategy 2 is an instance of ESP. For $s_0, s_1 > 0$ we define the smoothed count $s(x; x_{\leq k})$ of bit x and the smoothed total count t_k as follows

$$s(x; x_{\leq k}) := \begin{cases} \lambda s(x; x_{< k}) + 1, & \text{if } k > 0 \text{ and } x_k = x \\ \lambda s(x; x_{< k}), & \text{if } k > 0 \text{ and } x_k \neq x \\ s_x, & \text{if } k = 0 \end{cases} \quad \text{and } t_k := \begin{cases} \lambda t_{k-1} + 1, & \text{if } k > 0 \\ s_0 + s_1, & \text{if } k = 0 \end{cases}.$$

Strategy 2 predicts $p(x; x_{\leq k}) = s(x; x_{\leq k})/t_k$. In case $x_k = x$ we get

$$p(x; x_{\leq k}) = \frac{\lambda s(x; x_{< k}) + 1}{t_k} = \frac{\lambda t_{k-1}}{t_k} \frac{s(x; x_{< k})}{t_{k-1}} + \frac{1}{t_k} = \frac{t_k - 1}{t_k} p(x; x_{< k}) + \frac{1}{t_k},$$

similarly $p(x; x_{\leq k}) = \frac{t_k - 1}{t_k} p(x; x_{< k})$, if $x_k \neq x$. If we now choose $\alpha_k = \frac{t_k - 1}{t_k}$ and $p(x) = \frac{s_x}{s_0 + s_1}$ the above sequential probability assignment rule resembles (1). This insight allows us to adopt our analysis method. To do so, we require the following technical statement first.

Lemma 4.4. For $1 \leq a \leq b$ and $0 < \lambda < 1$ we have $\frac{1 - \lambda^a}{1 - \lambda^b} \geq \frac{a}{b}$.

Proof. Let $f(z) := \ln((1 - \lambda^z)/z)$, it suffices to prove that $f(a) \geq f(b)$. By $\ln \lambda^z \geq 1 - 1/\lambda^z$ we get $f'(z) = [(1 - \ln \lambda^z) \cdot \lambda^z - 1] / [a(1 - \lambda^a)] \leq 0$, so f is decreasing. \square

Corollary 4.5. Let \mathcal{S} be an arbitrary partition of $[1, n]$. Fix $0 < \lambda < 1$ and $m \geq 1$, define $t_k := \lambda t_{k-1} + 1$ for $k \geq 1$ and $t_0 = 1 + \lambda + \dots + \lambda^{m-1}$ for $k = 0$. If $\alpha_k = \frac{t_k - 1}{t_k}$ and Assumption 2.2 holds, then

$$\ell(x_{1:n}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:b}) \leq |\mathcal{S}| \cdot \left[\log \frac{n}{p(0)} + \frac{(\pi \log e)^2}{6 \log \frac{1}{\lambda}} + (n-1) \log \frac{1}{\lambda} \right]. \quad (15)$$

Proof. Let $k \geq 1$ and note that by $t_k = \lambda t_{k-1} + 1$ we may write $\alpha_k = \lambda t_{k-1}/t_k$ and $t_k = 1 + \lambda + \dots + \lambda^{k+m-1} = (1 - \lambda^{k+m})/(1 - \lambda)$ and get

$$\beta_i = \alpha_1 \cdot \dots \cdot \alpha_i = \frac{\lambda t_0}{t_1} \frac{\lambda t_1}{t_2} \dots \frac{\lambda t_{i-1}}{t_i} = \frac{t_0}{t_i} \cdot \lambda^i = \frac{1 - \lambda^m}{1 - \lambda^{m+i}} \cdot \lambda^i.$$

We now proceed by bounding the terms dependent on β_i in (6):

$$\beta_i = \frac{(1 - \lambda^m) \lambda^i}{1 - \lambda^{m+i}} \stackrel{\text{Lem. 4.4}}{\geq} \frac{m \lambda^i}{m+i} \stackrel{m \geq 1}{\geq} \frac{\lambda^i}{i+1} \quad \text{and} \quad \frac{\beta_i}{\beta_a} = \frac{1 - \lambda^{m+a}}{1 - \lambda^{m+i}} \cdot \lambda^{i-a} \stackrel{a \leq i}{\leq} \lambda^{i-a}$$

From the above inequalities we obtain

$$\log \frac{1}{\beta_{n-1}} \leq \log \frac{n}{\lambda^{n-1}} \quad \text{and} \quad \sum_{a < i < b} \log \frac{1}{1 - \beta_i/\beta_a} \leq \sum_{a < i < b} \log \frac{1}{1 - \lambda^{i-a}} \stackrel{\text{Lem. 4.1}}{\leq} \frac{(\pi \log e)^2}{6 \log \frac{1}{\lambda}}.$$

Finally we plug the above inequalities into (6) and rearranging yields (15). \square

For $k \rightarrow \infty$ we have $t_k \rightarrow \frac{1}{1-\lambda}$, thus $\alpha_k \rightarrow \lambda$, i. e. we expect the smoothed counts method to perform similar to ESP with fixed smoothing rate λ , when the input is large enough. Bound (15) reflects this behavior, it differs from (10) only by the additive term $|\mathcal{S}| \log n$. Furthermore, the optimal value of λ in (15) matches the optimal value of α in (10).

5 Experiments

For inputs of length n we experimentally checked the tightness of our bounds from the previous section for a wide range of ESP-instances with smoothing rate choices (i) fixed “optimal” smoothing rate $\alpha = \exp(-\pi/\sqrt{6(n-1)})$ (here “optimal” means that the corresponding bound, c. f. Corollary 4.2, is minimized), (ii) varying smoothing from Corollary 4.3 and (iii) varying smoothing from Corollary 4.5 with “optimal” $\lambda = \exp(-\pi/\sqrt{6(n-1)})$ and $m = 1$. Since our bounds from corollaries 4.2, 4.3 and 4.5 are worst-case bounds we compare them to the empirically measured (approximate) worst-case redundancy. Furthermore, we compare the (approximate) worst-case redundancy of (i), (ii) and (iii) to each other. We now explain the details below.

Experimental Setup. In the following let smoothing rate sequence $\alpha_{1:\infty}$, input length $n = 1000$, partition $\mathcal{S} = \{(0, 200], (200, 700], (700, 1000]\}$ and $\varepsilon = 0.05$ be fixed. (We inspected the outcome of our experiments for different parameters and got similar results, hence these values.) We want to judge on our bounds on a wide range of ESP-instances, in particular we choose class $\mathcal{C} = \{(\alpha_{1:\infty}, p) \mid 0 < \varepsilon \leq p(0), p(1)\}$ of ESP-instances. To do so, we have to modify our bound slightly, we must replace the term $p(0)$ by ε : For instance, in Situation (i), we may bound the redundancy of any $\text{ESP} \in \mathcal{C}$ of prefix $x_{1:k}$ of given $x_{1:n}$ as follows

$$\ell(x_{1:k}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:\min\{k,b\}}) \leq |\mathcal{S}| \cdot \left[\frac{2\pi \log e}{\sqrt{6}} \cdot \sqrt{k-1} + \log \frac{1}{\varepsilon} \right], \quad (16)$$

for $1 \leq k \leq n$. Since the resulting bounds remain worst-case bounds, we compare the resulting bounds for situations (i)-(iii) to the worst-case redundancy

$$r(k) := \max_{\text{ESP} \in \mathcal{C}, x_{1:n}} \left(\ell(x_{1:k}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:\min\{k,b\}}) \right). \quad (17)$$

Unfortunately, computing the maximum is intractable, since \mathcal{C} is uncountably infinite and there are exponentially many sequences $x_{1:n}$. To lift this limitation we take the maximum over a finite subset of ESP-instances from \mathcal{C} and inputs $x_{1:n}$, specified as follows: For numbers $q_0, \dots, q_{|\mathcal{S}|} \in \{0.05, \dots, 0.95\}$ we consider pairs $(\text{ESP}, x_{1:n})$ s. t. $\text{ESP}(0; \phi) = q_0$ (q_0 determines an ESP-instance) and $x_{1:n}$ is drawn uniform at random from all sequences where for the i -th segment $[a, b] \in \mathcal{S}$ subsequence $x_{a:b}$ has exactly $\lfloor q_i \cdot (b - a + 1) \rfloor$ 1-bits (q_i determines the (approximate) fraction of 1-bits in the i -th segment). We now take the maximum in (17) over all combinations $(q_0, \dots, q_{|\mathcal{S}|})$ and repeat the random experiment 100 times for every combination $(q_0, \dots, q_{|\mathcal{S}|})$ (in total $19^{|\mathcal{S}|+1} \cdot 100$ simulations). Figure 1 depicts the approximation of $r(k)$ (solid lines) and our bounds on $\ell(x_{\leq k}; \text{ESP}) - \sum_{[a,b] \in \mathcal{S}} h(x_{a:\min\{b,k\}})$ (dashed lines). (For instance, bound (16) is depicted as dashed line in the left plot of Figure 1.)

Approximate Worst-Case Redundancy. We now compare (approximate) $r(k)$ for smoothing rate choices (i)-(iii) and observe: On one hand, as long as k is small, varying smoothing rates, (ii) and (iii), yield lower redundancy than (i), and (iii) performs better than (ii). On the other hand, when k is large (i), (ii) and (iii) don’t differ too much. The increase in redundancy at $k = 201$ and $k = 701$ is nearly identical in all cases, the difference in redundancy is almost entirely caused by segment $(0, 200]$.

Bounds Behavior. Now we compare the bounds to (approximate) $r(k)$. In general, the tightness of our bounds decreases as the number of segments increases. This is plausible,

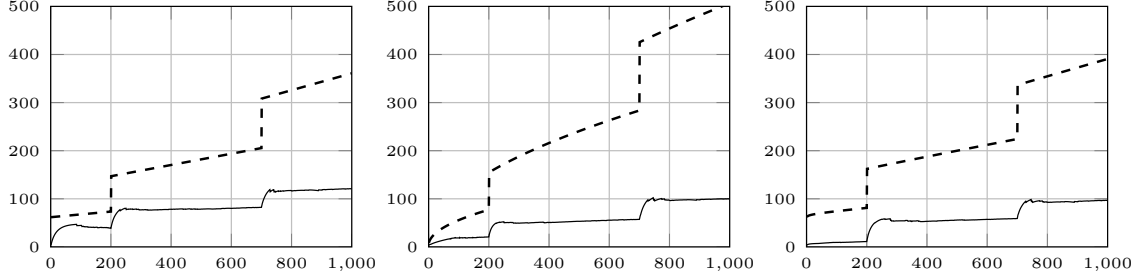


Figure 1: Redundancy bound (dashed line) and approximate worst-case redundancy $r(k)$ (solid line) of class $\{(\alpha_{1:\infty}, p) \mid 0 < \varepsilon \leq p(0), p(1)\}$ for $\varepsilon = 0.05$ w.r.t. competitor with partition $\mathcal{S} = \{[1, 200], (200, 700], (700, 1000]\}$ on the length- k prefix, $1 \leq k \leq n$ of a sequence with length $n = 1000$. The x -axis is prefix length k and the y -axis is redundancy in bit. Every plot corresponds to a different smoothing rate choice: (i) fixed “optimal” smoothing rate $\alpha = \exp(-\pi/\sqrt{6(n-1)})$, (ii) varying smoothing from Corollary 4.3 and (iii) varying smoothing from Corollary 4.5 with “optimal” $\lambda = \exp(-\pi/\sqrt{6(n-1)})$ and $m = 1$.

since we essentially concatenated the worst-case bound for $|\mathcal{S}| = 1$. However, we don’t know, whether or not the worst-case redundancy for $|\mathcal{S}| = 1$ can appear in multiple adjacent segments at the same time. Experiments indicate that this may not be the case. Furthermore, in (i) the bound is tightest, especially within segment $(0, 200]$. In cases (ii) and (iii) the bounds are more loose. An explanation is, that in the corresponding proofs we worked with rather generous simplifications, e. g. when bounding $-\sum_{a < i < b} \log(1 - \beta_i/\beta_a)$ from above. If we compare (ii) to (i) and to (iii) we can see, that bound (ii) is tighter for very small k . The reason is simple: Bound (ii) does not depend on a smoothing rate parameter, whereas (i) contains the term $1/\log \frac{1}{\alpha}$ and (iii) contains the term $1/\log \frac{1}{\lambda}$. These terms dominate the bounds, when k is small and α and λ are close to 1. (We have $\alpha = \lambda \approx 0.96$, since α and λ were chosen to minimize the corresponding bound for $n = 1000$.)

6 Conclusion

In this work we analyzed a class of practical and adaptive elementary models which assign probabilities by exponential smoothing, ESP. Our analysis is valid for a binary alphabet. By choosing smoothing rates appropriately our strategy generalizes count smoothing (Strategy 2) and probability smoothing from PAQ (Strategy 3). Due to its low memory footprint and linear per-sequence time complexity ESP is attractive from a practical point of view. From a theoretic point of view ESP is attractive as well: For various smoothing rate sequences ESP has redundancy only $O(s\sqrt{n})$ above any PWS with s segments, an improvement over previous approaches. A short experimental study supports our bounds.

Nevertheless, experiments indicate that there is room for an improved analysis. Despite minor technical issues a major approach would be to obtain redundancy bounds w. r. t. PWS that take the similarity of adjacent segments into account. That is, if adjacent segments have very similar distributions, the increase in redundancy should be small, compared to adjacent segments with drastically different distributions. Furthermore, it is desirable to generalize the analysis to a non-binary alphabet. We defer a thorough experimental study that compares ESP to other adaptive elementary models to future research.

Acknowledgement. The author thanks Martin Dietzfelbinger, Martin Aumüller and the anonymous reviewers for helpful comments and suggestions that improved this paper.

References

- [1] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2nd edition, 2006.
- [2] Paul G. Howard and Jeffrey S. Vitter. Analysis of arithmetic coding for data compression. In *Proc. Data Compression Conference*, volume 1, pages 3–12, 1991.
- [3] Paul G. Howard and Jeffrey S. Vitter. Practical Implementations of Arithmetic Coding. Technical report, Brown University, USA, 1991.
- [4] Christopher Mattern. Combining Non-stationary Prediction, Optimization and Mixing for Data Compression. In *Proc. International Conference on Data Compression, Communications and Processing*, volume 1, pages 29–37, 2011.
- [5] Christopher Mattern. On Probability Estimation via Relative Frequencies and Discount. 2013. <http://arxiv.org/abs/1311.1723>.
- [6] Eado Meron and Meir Feder. Finite-memory Universal Prediction of Individual Sequences. *IEEE Trans. on Information Theory*, 50:1506–1523, 2006.
- [7] Alexander O’Neill, Marcus Hutter, Wen Shao, and Peter Sunehag. Adaptive Context Tree Weighting. In *Proc. Data Compression Conference*, volume 22, pages 317–326, 2012.
- [8] Gil I. Shamir and Neri Merhav. Low-complexity sequential lossless coding for piecewise-stationary memoryless sources. *IEEE Trans. on Information Theory*, 45:1498–1519, 1999.
- [9] Joel Veness, Kee Siong Ng, Marcus Hutter, and Michael H. Bowling. Context Tree Switching. In *Proc. Data Compression Conference*, volume 22, pages 327–336, 2012.
- [10] Joel Veness, Martha White, Michael Bowling, and A. András Gyorgy. Partition tree weighting. In *Proc. Data Compression Conference*, volume 23, pages 321–330, 2013.
- [11] Frans Willems, Yuri M. Shtarkov, and Tjalling J. Tjalkens. The context-tree weighting method: basic properties. *IEEE Trans. on Information Theory*, 41:653–664, 1995.
- [12] Frans M. J. Willems. Coding for a binary independent piecewise-identically-distributed source. *IEEE Trans. on Information Theory*, 42:2210–2217, 1996.